

# A abstração de ordenação causal e uma maneira simples de implementá-la

Michel Rayrai

*IRt,\$A, Beaulieu Campus, F-35042 Rennes Cedex, França*

Outros Schippers

*Escola Politécnica F•d•e!c, Departamento de InJomiatiqve, C ff -1015 Lausanne, França*

Sam Toweg

*Departamento de Ciência da Computação, Cornell University, Três, NT 94853, ttSA*

Comunicado por L. Kott Recebido

em 11 de maio de 1990

Revisado em 3 de abril de 1991

## resumo

Raynal, M., A. Schiper e S. Toueg, A abstração de ordem causal e uma maneira simples de implementá-la, *Information .roc\*.ssing \* •.tters* 39 (1991) 343-350.

O controle em sistemas distribuídos é introduzido principalmente para reduzir o não determinismo. Este não determinismo deve-se, por um lado, à execução assíncrona dos processos localizados nos vários sites do sistema e, por outro lado, à natureza assíncrona dos canais de comunicação. Para limitar o assincronismo devido aos canais de comunicação, uma nova relação de ordenação de mensagens, conhecida como ordenação causal, foi introduzida por Birman e Joseph. Depois de dar alguns exemplos dessa ordenação causal, propomos um algoritmo simples para implementá-la. Este algoritmo é baseado na numeração da sequência de mensagens. Uma prova da correção do algoritmo também é fornecida.

*Palavras-chave:* Computação distribuída, algariiii'ur• distribuído, eausai or3.•âog, n•-twork produ!

## 1. Controlando o não detenninismo

Controlar o não determinismo é uma das tarefas essenciais em sistemas computacionais. Em sistemas centralizados o não determinismo é causado pelas interrupções. Em sistemas distribuídos o não terminismo também se deve à execução assíncrona dos processos paralelos juntamente com a natureza assíncrona dos canais de comunicação. Além disso, tais sistemas também podem ser afetados por ocorrências não determinísticas de falhas de componentes.

Para reduzir o assincronismo dos processos e resolver os problemas relacionados (por exemplo, problema de alocação de recursos), geralmente se introduz a sincronização. mecanismos [14,12]. Para reduzir o assincronismo dos canais de comunicação é possível construir sincronizadores de rede que forcem processos e canais a progredir em passos sincronizados [1,7]. Isso corresponde a definir uma máquina virtual distribuída na qual o comportamento indesejável devido ao não terminismo foi suprimido. O mesmo vale para os protocolos

efetivamente para mascarar canais não confiáveis: os protocolos de comunicação definem uma abstração de "canal confiável" que é implementada sobre um canal não confiável. O protocolo de bit alternativo [2] ou o protocolo de Stenning [16] são alguns exemplos: ambos constroem um canal sem perda, duplicação ou desequenciamento de mensagem em cima de um canal que pode perder ou duplicar mensagens, ou ainda, no caso do protocolo de Stenning, que pode desequenciá-los. A abstração confiável do canal, assim como os mecanismos de sincronização, elimina os comportamentos indesejáveis.

Neste artigo, consideramos um sistema distribuído confiável e estamos preocupados com a realização de um esquema de comunicação que elimine um determinado comportamento indesejável. Mais precisamente, estamos preocupados com a abstração chamada "ordenação causal", conforme proposto pelo sistema Isis [3]. O artigo está organizado da seguinte maneira. Na Seção 2, definimos ordenação causal e mostramos sua utilidade. Na Seção 3 mostramos uma maneira fácil e natural de implementar a ordenação causal (que difere da implementação de Isis), e damos na Seção 4 uma prova de sua correção.

## 2. Ordenação causal

### 2.1. Definição

A ordenação causal de eventos em um sistema distribuído (um evento correspondente à emissão de uma mensagem, a recepção de uma mensagem ou uma ação interna) é baseada na conhecida relação "aconteceu antes", notada  $\rightarrow$  [9]. Se  $A$  e  $B$  são dois eventos, então  $A \rightarrow B$  se e somente se uma das seguintes

condições for verdadeira: (i)  $A$  e  $B$  são dois eventos ocorrendo no mesmo

local,  $A$  ocorre antes de  $B$ ; (ii)  $A$  é a emissão de uma mensagem, e  $B$  corresponde à recepção da mesma mensagem; (iii) existe um evento  $C$  tal que  $A \rightarrow C$  e  $C \rightarrow B$ .

A ordenação causal diz respeito à ordem na qual duas mensagens  $M_1$  e  $M_2$  são entregues. Vamos introduzir a notação **DELIVERY**( $M_i$ ) para indicar o evento correspondente à entrega de  $M_i$ .

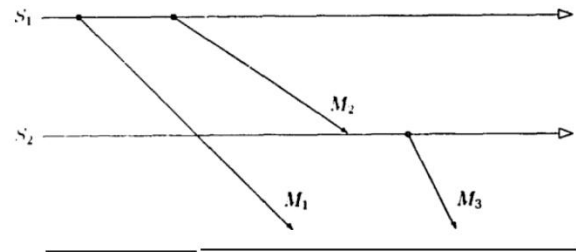


Fig. 1. A ordenação causal é respeitada se  $M_1$  for entregue antes de  $M_2$ .

ao seu processo de destino, e considere dois eventos **SEND**( $M_1$ ) e **SEND**( $M_2$ ). A ordem causal é respeitada se: se

**SEND**( $M_1$ )  $\rightarrow$  **SEND**( $M_2$ ), e  $M_1$  e  $M_2$  tem mesmo destino, então **DELIVERY**( $M_1$ ) **ENTREGA** ( $M_2$ )

Em outras palavras, no caso de uma relação "aconteceu antes" entre duas emissões para um mesmo destino, existe uma relação "aconteceu antes" entre a entrega das mensagens. Por exemplo, na Fig. 1, a ordem causal é respeitada se a mensagem  $M_1$  for entregue no site S3 antes da mensagem  $M_2$ .

Ao enviar mensagens usando protocolos de transporte tradicionais, a ordem causal pode ser claramente violada. Implementar ordenação causal significa adicionar um protocolo ao sistema original de forma que a ordenação causal nunca seja violada no sistema equipado com o protocolo sobreposto. Para evitar confusão entre os dois sistemas, usaremos a seguir consistentemente a expressão mensagem "recepção" para denotar uma mensagem recebida pelo sistema original, e mensagem "entrega" para denotar uma mensagem recebida pelo sistema equipado com o protocolo de ordenação causal. Uma analogia pode ser feita aqui entre canais FIFO e um sistema implementando ordenação causal. Um canal **FIFO** reduz o não determinismo de uma comunicação ponto a ponto. Um sistema implementando ordenação causal globalmente reduz a indeterminação das comunicações.

### 2.2. Exemplos da utilidade da ordenação causal

#### Gerenciamento de dados replicados

Considere um dado  $N$  replicado em vários locais. A fim de assegurar a consistência mútua de  $N$  vari-

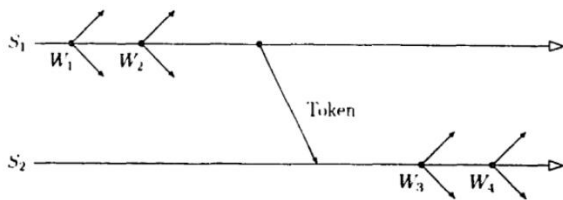


Fig. 2. Gestão de dados replicados.

ou cópias  $m_1, \dots, m_n$ , as atualizações dessas cópias devem ocorrer na mesma ordem. Em particular, isso introduz a necessidade de exclusão mútua nas atualizações feitas por dois sites diferentes. Isso pode ser resolvido usando um token; quando de posse do token, um site pode atualizar sua cópia e transmitir a atualização para todos os outros sites que tenham uma cópia de  $A$  (Fig. 2). A ordenação causal garante que todos os sites recebam todas as atualizações na mesma ordem (IN, wi-fi não será entregue antes de  $l_i$ , para  $i < j$ ), o que garante a consistência mútua das cópias [4,8]. O token garante a ordenação total do update e a ordenação causal garante que todas as cópias  $x_i$  sejam atualizadas nesta mesma ordem.

#### Observando um sistema distribuído

A ordenação causal também pode ser usada para fornecer observações consistentes de um sistema distribuído.

Considere vários sites monitorando eventos que ocorrem localmente e suponha que todas essas informações sejam de interesse de um sistema operacional de site observador.

Se as informações de monitoramento coletadas pelos sites forem enviadas ao OS respeitando a ordem causal, as informações serão recebidas em uma origem coerente (ou seja, em uma ordem que não viole a ordem de precedência dos eventos). Por exemplo, na Fig. 1, o observador  $S$  recebe mensagens de sondagem  $M$  e depois  $O_1$ , ou seja, em uma ordem de observação consistente.

#### Alocação de recursos

Considere o problema de exclusão mútua ou, mais geralmente, um problema de alocação de recursos. Isso pode ser resolvido por um alocador de recursos residente em um site, que recebe solicitações de alocação em uma fila **FIFO**. Com esse esquema, os pedidos são atendidos na ordem de recebimento. Isso pode, no entanto, ser considerado insatisfatório em alguns casos: pode ser desejável honrar os pedidos não pela ordem de recepção, mas pela ordem de emissão, ou seja, se dois

$A_i$  deve e  $R$  são tais que  $R_i \rightarrow A_i$ , então solicita ser honrado antes de  $A_i$ . Este é novamente um problema de ordenação causal (considere na Fig. 1 que  $M_1$  carrega a solicitação  $A_1$  e  $M_2$  carrega a solicitação  $A_2$ ).

Observe que este exemplo, no caso particular de exclusão mútua, foi usado por Lamport em seu artigo introduzindo relógios lógicos [9]. A solução de Lamport foi ordenar totalmente as solicitações, o que não é necessário se as mensagens forem ordenadas causalmente.

### 3. Implementação de ordenação causal

#### 3.1. Trabalho relatado

A implementação de Isis da ordenação causal [3] é direta: uma mensagem carrega consigo toda a história das comunicações que a precederam.

Considere novamente a Fig. 1. A mensagem  $M$  enviada por  $S_2$  para  $S$  carrega a seguinte informação sobre seu passado: mensagem  $M_1$ , com a informação de que foi enviada para  $S_3$ , e mensagem  $M_2$  com a informação de que foi enviada para  $S_2$ . Ao receber  $M_3$ , o Site  $S_3$  também recebe uma cópia de  $M$ : se ainda não foi entregue,  $M$  será entregue naquele momento, ou seja, antes de  $M$ . Para evitar o crescimento ilimitado das informações adicionadas a uma mensagem, deve ser adicionado um mecanismo para informar as mensagens que foram recebidas e entregues: essas mensagens não precisam mais ser enviadas. Este mecanismo é, no entanto, bastante complexo. De qualquer forma, o tamanho do pneu das informações adicionadas às mensagens é ilimitado.

Outra implementação de ordenação causal é dada em [15]. Nesta implementação, o controle na formação (e não nas mensagens como no Isis) é adicionado às mensagens. Esta informação de controle permite ao local de destino de uma mensagem não saber se existem mensagens que têm de ser entregues antes de  $l_i$ , de forma a respeitar a ordem causal: se for o caso, não é imediatamente entregue. A informação de controle é composta por um número limitado de pares (local de destino, vetor tempo), onde tempo é um tempo que define uma ordenação parcial dos eventos de um sistema distribuído [5,11]. A implementação proposta neste artigo tem algumas semelhanças com esta implementação. Não usamos tempos vetoriais e, portanto, esse algoritmo é muito mais fácil de entender.

### 3.2. Uma implementação simples

Considere que temos um sistema não causal original (isto é, onde a ordem causal pode ser violada) NC. Nosso objetivo é definir um protocolo no sistema NC para construir um novo sistema C no qual a ordenação causal nunca seja violada. Para implementar o protocolo, damos a cada site em informações de controle NC representando a percepção do site do estado do sistema (mais precisamente, a percepção do site das comunicações). Esta informação permitirá que cada site decida quando uma mensagem recebida pode ser entregue ao sistema C. Além disso, as informações adicionadas a cada mensagem permitirão que o site receptor no NC, quando uma mensagem for entregue, atualize sua percepção do estado do sistema.

#### Informações locais de um

site Cada site no sistema NC gerencia as duas variáveis seguintes (onde  $ii$  é igual ao  $r$ : número de sites no sistema):

**DELIV**: array[1..n] de inteiro;

(' inicialmente  $DELIV[i] = 0$  para todo  $i$  \* )

**SENT**: array[1..n, 1..n] de inteiro;

(' inicialmente  $SENT[i, y] = 0$  para todo  $i, y$  \* )

Usaremos a notação **DELIV**, e **SENT**, para nos referirmos às variáveis gerenciadas pelo site  $S_i$ . Então no site  $S_i$ , a variável  $DELIV_i[y]$  representa o número de mensagens enviadas de  $S_i$  e entregues a  $S_y$ , a variável  $SENT_i[k, l]$  representa o conhecimento de  $S_i$  do número de mensagens enviadas (mas não necessariamente entregue) de  $S_k$  para  $S_l$  (mais precisamente  $S_k$  foi informado que pelo menos  $SENT_i[k, l]$  mensagens  $z'$  ere enviadas de  $S_k$  para  $S_l$ ).

Para ilustrar essas definições, considere a Fig. 3 e, mais precisamente,  $SENT_i[3, 1]$  que é igual a 1, significando que o site  $y$  sabe que uma mensagem (aqui fill) foi enviada do site  $k = 3$  para o site  $i = 1$ . Como veremos, esta informação foi enviada junto com  $\tilde{n}f2$ .

#### Comportamento de

um site O comportamento de um site  $S$  é expresso pelas duas regras a seguir, que regem a emissão e a recepção de uma mensagem no sistema NC.

(“i) Emissão de uma mensagem  $M$  de  $S_i$  para  $S_j$ :

enviar(  $Se, ENVIADO_i$  ) para  $S_j$ ;

$ENVIADO_i[i, y] := ENVIADO_i[i, y] + 1$ ;

Assim,  $M$  é enviado junto com uma cópia de  $SENT_i$ , ou seja, com o site  $S_j$ 's visualização do estado do sistema.

(“u) Recepção de  $M$ ,  $ST f$  ) enviado de  $S$  para  $S_j$ , onde  $STD$  representa a informação de controle “SENT” transportada pela mensagem  $M$  (ver Fig. 3):

espera(para todo  $k, DELIV_i[k] \geq STD[\tilde{n}, i]$ );

entrega de  $M$  ao sistema C;

$DELIV_i[y] := DELIV_i[y] + 1$ ;

$ENVIADO_i[j, f] := ENVIADO_i[j, f] + 1$ ;

para todo  $k, 1: SENT_i[k, l]$

$= \max(ENVIADO_i[k, l], STD[k, f])$ ;

Assim, uma mensagem  $\tilde{n}f$  enviada por  $f_i$ , pode ser entregue a  $S_j$  somente se para todo  $k, DELIV_i[k] \geq STD[k, j]$ , que expressa que todas as mensagens causalmente precedentes, cuja existência é revelada por  $STD$ , foram entregue. Considere, por exemplo, a mensagem  $M$  na Fig. 3.  $STp_i[3, 1] = 1$ , o que significa que  $M$  pode ser entregue ao site 1 quando  $DELIV_i[3] \geq 1$ . Isso é verdade assim que  $\tilde{n}f4$  é entregue. Observe neste exemplo particular que  $DELIV_i[3] > STD_i[3, 1]$  quando  $\tilde{n}f4$  é recebido no site 1. Esta situação se deve à existência de  $3f3$ , cuja emissão não precede causalmente a emissão de  $M$ :  $SEND(\tilde{n}f3) \rightarrow ENVIAR(\tilde{n}f4)$ . Os dois eventos **SEND**( $\tilde{n}f3$ ) e **SEND**( $\tilde{n}f4$ ) são considerados con-

atual.

Depois que uma mensagem é entregue, a percepção do site sobre o estado do sistema é atualizada de maneira óbvia.

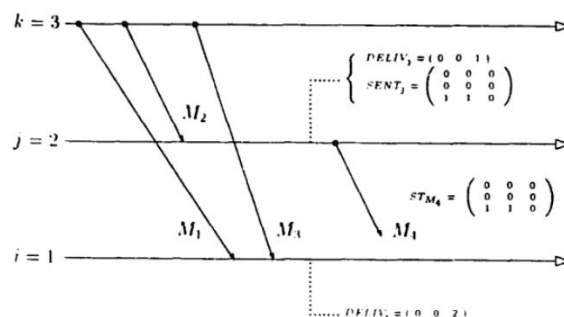


Fig. 3. Ilustração das estruturas de dados.

#### 4. Prova de correção

A correção do algoritmo acima será provada nas duas etapas usuais; segurança e vivacidade. No entanto, antes de desenvolver essas provas, precisamos de alguns resultados intermediários que serão dados na Seção 4.1. Em seguida, provaremos na Seção 4.2 a segurança de nosso protocolo, o que corresponde a provar que a ordenação causal nunca é violada. Finalmente na Seção 4.3 provaremos a vivacidade do algoritmo, que corresponde a provar que toda mensagem será eventualmente entregue. Com relação ao sistema NC subjacente no qual a ordenação causal é implementada, precisamos assumir que todas as mensagens enviadas são recebidas, ou seja, nenhuma mensagem é perdida no NC e o sistema NC está ativo.

##### 4.1. Resultados preliminares

A fim de provar a segurança e vivacidade do nosso protocolo, precisamos de alguns resultados preliminares.

**Proposição 1.** Considere uma mensagem  $M$  sent por  $S$  para  $S_j$ ; tal que  $\text{STD}[k, i] > 0$  (incluindo  $k \rightarrow j$ ), e seja  $\text{ST} \gg k, i \rightarrow x$ . Então  $M$  é a  $x$ -ésima mensagem enviada por  $S_k$  para  $S_j$ , é tal que  $\text{SEND}(M_f) \rightarrow \text{SEND}(l_f)$ .

**Prova.** (i)  $k \rightarrow j$ . De  $\text{STp}[j, i] = z$  e do protocolo de emissão,  $\text{SENT}(y, i) = z$  no momento em que  $M$  é enviado. Além disso, segue diretamente do protocolo que  $\text{SENT}(y, i)$  nunca é incrementado ao receber uma mensagem por  $S_j$ . Assim, no momento em que  $\tilde{n}$  é enviado, o protocolo de emissão foi executado  $x$  vezes por  $S$ , ou seja,  $M$  foi enviado antes de  $M$ , e  $\text{SEND}(3f) \rightarrow \text{SEND}(M)$ .

(ii)  $k \neq j$ . A prova consistirá na construção de uma cadeia causal  $\text{SEND}(3f) \rightarrow \text{SEND}(l_f)$ .

Seja  $N_j$  a primeira mensagem enviada por  $S_j$ ; tal que  $\text{IT } g[k, i] \rightarrow x$  (Fig. 4), ou seja,  $\text{SEND}[k, i] \rightarrow x$  quando  $A$  foi enviado ( $A$ , pode ser  $3f$ ). Assim, antes das emissões de  $A$ ,  $S$  deve ter recebido uma mensagem  $My$  de  $S$  com  $\text{STD}[k, i - x]$  (consulte a última etapa do protocolo de entrega) e  $\text{SEND}(\tilde{n}, )$

ENVIAR( $3f$ ).

Porque  $\text{STD}[k, i - \text{STD}[k, i]]$ , a mesma construção aplicável a  $l_f$  pode ser aplicada a  $My$ , e assim por diante. Cada iteração desta construção introduz um novo site  $S_t(s)$  diferente de todos os anteriores

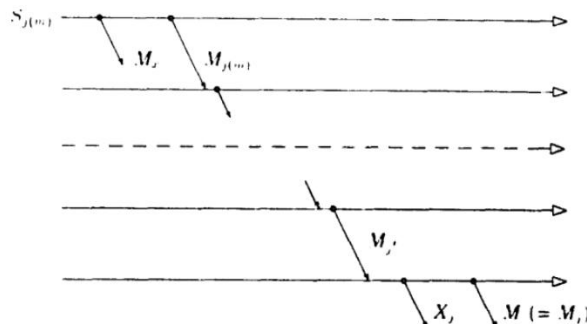


Fig. 4. Construção da cadeia causal  $\text{SEND}(M_f) \rightarrow \text{SEND}(l_f)$ .

sites, pois  $\bullet$  define uma ordem [9]: se  $S_t, S_j, i/2$  for  $s_1 s_2$ , teríamos um ciclo na cadeia causal construída.

Como o número de sítios é finito, finalmente terminamos com  $S_{(m)} = S$ ,  $\text{IT} \bullet$  enviado por  $S(q, e \text{ STMG}(k, i \rightarrow x$ , que foi tratado em (i). Assim, a cadeia causal  $\text{SEND}(\tilde{n}, )$

$\text{SEND}(f_j, p \rightarrow) \rightarrow \text{SEND}(\tilde{n}, ) \rightarrow \text{SEND}(M)$  foi criado. O

Os dois lemas a seguir são resultados intermediários necessários para provar a Proposição 4.

**Lema 2.** Seja  $M$  a  $x$ -ésima mensagem enviada por  $S_t$  para  $S$ ,  $x > 0$ , e considere uma mensagem  $M$  tal que  $\text{SEND}(3fg) \rightarrow \text{SEND}(\tilde{n}, )$ . Então  $\text{STD}[j, i] > x$ .

**Prova** (esboço). Do protocolo de emissão, segue imediatamente que  $\text{STD}[i, y] = z - 1$ . A partir disso, o lema pode ser facilmente comprovado considerando o protocolo de emissão e recepção aplicado a qualquer conexão de cadeia causal  $\bullet \text{ SF.ND}(M, )$  Em  $\text{SEND}(\tilde{n}, )$ . O

**Lema 3.** Considere uma mensagem  $M$  e  $\text{ef STD}[j, i] \rightarrow x$ . Então  $M$  a mensagem  $\{x + 1\}$  enviada por  $S$ , para  $S$ , é tal que  $\text{SEND}(3fg+) \rightarrow \text{SEND}(\tilde{n}, )$ .

**Prova.** Suponha que  $\text{SEND}(3fy + t) \rightarrow \text{ey} \text{ SEND}(M)$ .

Lema 2 temos  $\text{ST3f}[i, i] > px + 1$ , o que mostra a contradição. O

**Proposição 4.** Considere  $M$  enviado por  $S$  para  $S_j$ , e outro sítio  $S_z$  (incluindo  $z \rightarrow k \rightarrow j$ ). Então  $\text{STD}[k, f]$  é igual ao número de mensagens  $N$  enviadas por  $S_j$  para  $S$ , tal que  $\text{SEND}(3f) \rightarrow \text{SEND}(\tilde{n}, )$ .

Ptmf. Seja  $m$  o número de mensagens  $M'$  enviadas por  $S_z$  para  $S$ , tal que  $\text{SEND}(C_f) \rightarrow \text{SEND}(\tilde{n}f)$ . e seja  $\text{STE } k, i - x$ . (i)

Pela Proposição 1, o . a mensagem enviada de  $S$  para  $S$ , é tal que  $\text{SEND}(\tilde{n}f, ) \rightarrow \text{SENDt}C_f$ . Assim  $m \geq z$ . (ii) Considere  $M$

de  $S$ , para  $S$ , a  $(x + 1)^{\text{a}}$  mensagem enviada Pelo Lema 3.  $\text{SENO}(\tilde{n}f, qt) \rightarrow \text{SEND}(3f)$ . Assim  $m < x$ .

De (i) e (ii) concluímos que  $n_i = z$ . O

Corolário 5. Considere um arquivo de mensagens se i front  $St$  para  $S$ , e uma mensagem  $M$ ; tal que  $\text{SENO}(\tilde{n}f, ) \rightarrow \text{SENO}(\tilde{n}f)$ .  $\forall i \in \text{STp}[y, i] < \text{STD}[y, eu]$ .

**Prmf.** Pela Proposição 4.  $\text{STD}[y, i]$  é igual ao número de mensagens  $M'$  enviadas de  $St$  para  $S$ , tal que  $\text{SENO}(\tilde{n}f) \rightarrow \text{SENO}(\tilde{n}f, )$ . Mas temos  $\text{SEND}(\tilde{n}f) \rightarrow \text{SEND}(3f) \rightarrow \text{SENDt fill}$ . Assim, porque  $\tilde{n}f$  também é enviado de  $S$ , para  $S$ , pela Proposição 4,  $\text{STp}[y, i]$  é pelo menos igual a  $\text{STD}[y, i] + 1$ . Isso permite concluir que  $\text{ST}[y, i] < \text{STp}[y, i]$ . O

Proposição 6. Considere um itinessage  $M$  enviado de  $St$  para  $S$ .  $IJ$  nem  $M$  nem qualquer mensagem  $M'$  tal que  $\text{SENO}(M) \rightarrow \text{SENDt}M'$  foi entregue a  $S$ , 'depois  $\text{DELIV}[y] \text{ y } \text{ST}[y, i]$ .

**Prova.** Temos apenas que considerar as mensagens  $3f$  enviadas por  $St$ , apenas a entrega dessas mensagens pode aumentar o  $\text{DELIV}[y]$ . Por Proposição 4,  $\text{STp}[y, i]$  é igual ao número de mensagens  $\tilde{n}f$  enviadas por  $S$ , a fi, tais que  $\text{SFV If } M'' \rightarrow \text{SEND}(3f)$ . Por hipótese, somente essas mensagens  $\tilde{n}f$  poderiam ter sido entregues a  $S$ . Assim  $\text{DELIV}[y] \leq \text{STD}[y, i]$ . O

Finalmente, a Proposição 7 será usada apenas para provar a vivacidade do nosso protocolo. Por esta razão podemos, na Proposição 7, assumir que nosso protocolo é seguro, ou seja, se  $\text{SEND}(3f) \rightarrow \text{SEND}(\tilde{n}f2)$  anJ  $MM$ , são enviados para o mesmo destino, então  $\tilde{n}f2$  Não pode ser entregue antes de  $My$  (safety será provado na Seção 4.ii).

**Proposição 7.** Suponha que nosso protcCol seja seguro, e considere .  $W$ , a  $x$ -ésima mensagem enviada de  $S$  para  $S$ ,

( $z > 0$ ). Então se  $\text{DELIV}[y] < x$ .  $M$  não foi entregue a  $S$ .

Prova.  $\text{DELIV}[y]$  é incrementado cada vez que uma mensagem de  $Si$  é entregue. Assim, se  $\text{DELIV}[y] < x$ , menos que  $x$  mensagens de  $Si$  foram entregues a  $S$ . Como o protocolo é considerado seguro, no máximo  $(x - 1)$  as primeiras mensagens de  $Si$  poderiam ter sido entregues a  $S$ . Assim,  $\tilde{n}f$ , não foi entregue a  $S$ . O

#### 4.2. Segurança

A segurança corresponde a provar que a ordem causal nunca é violada. isto é, se  $\text{SEND}(\tilde{n}f, ) \rightarrow \text{SENO}(\tilde{n}f2)$  e arquivo  $M$ ; são enviados para o mesmo destino. então  $\tilde{n}f2$  não pode ser entregue antes de  $M$ . Isso pode ser reformulado da seguinte maneira. Considere uma mensagem  $M$  enviada de  $S$ , para  $S$ , e seja  $M$  denotar a enésima mensagem entregue a  $S$ . Então, para todo  $i$ , enquanto  $\tilde{n}f$  não for entregue,  $\text{SENO}(\tilde{n}f) \rightarrow \text{SENO}(\tilde{n}f)$  não pode ser verdadeiro. Esta fórmula é adaptada para uma prova por indução.

(i) **Base stef** ( $r_i = 1$ ). Vamos supor que  $\tilde{n}f$ , a primeira mensagem entregue a  $S$ , é tal que  $\text{SENO}(\tilde{n}f) \rightarrow \text{SEND}(\tilde{n}f, )$ . e mostre que chegamos a uma contradição. Pela Proposição 4, como  $M$  é enviado de  $S$  para  $S$ ,  $\text{STp}[y, i] = 1$ . Inicialmente  $\text{DELIV}[y] = 0$ , então quando  $M$  é recebido por  $S$ ,  $\text{DELIV}[y] < \text{STD}[y, i]$ : a condição de entrega (Seção 3.2) não é satisfeita, ou seja,  $M$  não pode ser entregue, o que mostra a contradição.

(ii) **Etapa de indução.** Por hipótese de indução, enquanto  $\tilde{n}f$  não for entregue, para todo  $pn$ ,  $\text{SEND}(\tilde{n}f) \rightarrow \text{SENO}(\tilde{n}f)$  não pode ser verdadeiro. Vamos provar que o mesmo vale para  $ii + 1$ . Suponhamos novamente que  $\text{SEND}(3f) \rightarrow \text{SEND}(\tilde{n}f, )$  e mostre a contradição.

Antes de  $\tilde{n}f$ ,  $+ t$  ser entregue por  $S$ , temos pela Proposição 6:  $\text{DELIV}[y] \leq \text{STD}[y, i]$ . Ainda mais pelo Corolário 5 temos  $\text{STp}[y, i] < \text{STp}[y, i]$ , ou seja,  $\text{DELIV}[y] < \text{STp}[y, i]$  antes da recepção de  $\tilde{n}f + t$ , por  $S$ . Assim, a condição de entrega não é satisfeita (Seção 3.2), ou seja,  $\tilde{n}f + t$ , não pode ser entregue. Isso novamente mostra a contradição.

### 4.3. *hiueness*

Como já dito, provaremos aqui que toda mensagem é eventualmente entregue. Considere duas mensagens  $M'$  e  $M''$  e defina  $3f' < M''$  iff  $\text{SEND}(\tilde{n}f') \rightarrow \text{SEND}(\tilde{n}f'')$ .

Suponhamos que uma mensagem seja entregue assim que for recebida e sua condição de entrega se torne verdadeira. Portanto, uma mensagem não pode ser entregue se não for recebida ou se for recebida e sua condição de entrega for falsa.

Considere então um site  $S$ , e todas as mensagens enviadas para ele que não podem ser entregues. Usando a relação  $<$ , seja  $M$  uma das menores dessas mensagens (pode haver mais de uma dessas mensagens, pois  $<$  é uma ordem parcial). Esta mensagem  $3f$  foi enviada por algum  $St$ . Quando recebida, se  $3f$ , não puder ser entregue, temos (consulte a

Seção 3.2):

Reconhecimento

Os autores gostariam de agradecer aos revisores por suas sugestões que ajudaram a melhorar a clareza do artigo.

Referências

[1] B. Awerbuch, Complexidade da sincronização de redes, *J. ACM* 32 (4) (1985) 804-823.

[2] KA Barlett, RA Scantlebury e PT Wilkinson. Uma observação sobre transmissão full duplex confiável em links half duplex, *Comm. ACM* 12 (5) (1969) 260-261.

[3] K. Birman e T. Joseph, Comunicações reabilitadas na presença de falhas, *ACM Trans. Comput. ciência* 5 (1) (1987) 47-76.

[4] K. Birman e T. Joseph, Explorando a replicação em sistemas distribuídos, em: S. Mullender, ed., *Distributed Systems* (ACM, Nova York, 1990).

[5] C. Fidge, Timestamps em sistemas de troca de mensagens que preservam a ordenação parcial, em: *Prac. 11th Australian Computer Science Conf.*, University of Queensland, 1988.

[6] JM Helary, C. Jard, N. Plouzeau e M. Raynal, Detecção de propriedades estáveis em aplicações distribuídas, em: *Prac. 6º ACM Simp. em PODC* (1987) 125-136.

[7] JM Helary e M. Raynal, *Sincronização e Controle de Sistemas e Programas Distribuídos* (Wiley, Nova York, 1990), p. 200.

[8] T. Joseph e K. Birman, gerenciamento de baixo custo de dados replicados em sistemas distribuídos tolerantes a falhas, *ACM Trans. Comyul. ciência* 4 (1) (1986) 54-70.

[9] L. Lamport, Tempo, relógios e ordenação de eventos em um sistema distribuído, *Comm. ACM* 21 (7) (1978) 5g-ssh.

[10] F. Assuntos, Algoritmos para detecção de terminação distribuída, *Computação distribuída*. 2 (3) (1987) 161-175.

$3k$  tal que  $\text{DELIV}_i[k] < \text{STE}_i(k, i)$ .

Seja  $\text{STE}_i(k, i) = x$ , e considere  $M$ , a  $x$ -ésima mensagem enviada de  $S_z$  para  $S_i$ . De  $\text{DELIV}_i[k] < x$  e da Proposição 7 deduzimos que  $M$ , não foi entregue a  $S_i$ . Além disso, pela Proposição 1, temos  $\text{SEND}(3f_i) \rightarrow \text{SEND}(\tilde{n}f_i)$ , ou seja,  $Mg < M$ .

Isso significa que  $i$  hat  $\tilde{n}f$  não é uma das menores mensagens que não podem ser entregues a  $S_i$ , o que mostra a contradição. Segue-se que toda mensagem recebida por  $S$  eventualmente será entregue.

### 5. CONCLUSÃO

A ordenação causal das mensagens reduz o não determinismo decorrente da asiricronia dos canais de comunicação. Alguns exemplos foram dados para ilustrar a utilidade dessa noção de ordenação. Um algoritmo simples implementando esta ne ção, baseado na contagem de mensagens, foi proposto. Com essa implementação, o site de destino de uma mensagem sabe, ao receber uma mensagem, se ela pode ser entregue imediatamente ou se as mensagens precedentes ainda estão em andamento. Em contraste com a implementação potencialmente ilimitada do Ísis, a informação adicionada às mensagens para garantir a ordenação causal é aqui uma matriz  $n \times n$ , onde  $ii$  é o número de sites (cada entrada da matriz, no entanto, não é limitada). •. mesmo limite

- [11] F. Mattem, Tempo e locais globais de sistemas distribuídos, em:  
*Proc. Internat. Workshop on Parallel and Distributed Algorithms*,  
Bonas, France, 1988 (North-Holland, Amsterdam) 215-226.
- [12] G. Neiger e S. Toueg, Substituindo conhecimento *comum*  
e em tempo real em sistemas distribuídos, em: *Proc. 6º*  
**ACM** *Simp. em PODC* (1987) 281-293.
- [13] M. Raynal, *3 Algorithm for Mutual Exclusion* (MIT Press,  
Cambridge, MA, 1986), p. 107.
- [14] M. Raynal, *Disributed Computation and Metwor-ks:*  
*Concepts , Tools and Algorithms* (MIT Press, Cambriclge,  
MA, 1958), p. 166.
- [15] A. Schiper, J. Eggii e A. Sandoz, Um novo algoritmo para  
implementar a ordenação causal, em: *Foror. 3º Interno.*  
*Workshop on Distributed Algorithms*, Nice, **Lecture** Notes in Computer  
Science **392** (Springer, Berlin, 1989) 219—232.
- [16] W. Stunning, Um protocolo de transferência de dados, *Computer*  
*Networks* 1 (1976) 99-110.